



SDC 

STORAGE DEVELOPER CONFERENCE

SNIA  SANTA CLARA, 2017

“Performance bugs” and How to Find Them

Rune Erlend Jensen

MemoScale / Norwegian University of Science and Technology

Per Simonsen

MemoScale

Abstract

- ❑ “Performance bugs” are performance anomalies.
- ❑ We identify memory address aliasing as an important mechanism causing performance bias on modern Intel microarchitectures.
- ❑ We link this effect to memory allocators and demonstrate that they can hurt performance in several cases.



Presentation Overview

- ❑ Who are we?
- ❑ Optimization and data storage
- ❑ Optimization results
- ❑ Finding a "Performance Bug"
- ❑ Exploring the impacts



Who are the presenters?

- ❑ Rune Erlend Jensen, Performance engineer
 - ❑ Extensive optimization experience
 - ❑ Ph.D. candidate at the Norwegian University of Science and Technology.
- ❑ Per Simonsen, CEO, MemoScale
 - ❑ MSc Entrepreneurship/Computer Science
 - ❑ MSc Business Administration



MemoScale

- ❑ Start-up with strong ties to the Norwegian University of Science and Technology.
- ❑ *We leverage our expertise in optimization, software engineering and coding theory to solve bottlenecks and improve performance in data storage systems.*



Improving Performance

- ❑ Fastest unbroken SHA3 Candidate: Blue Midnight Wish
- ❑ Fastest Reed Solomon Erasure Coding implementation
- ❑ New erasure codes optimized for fast compute and rebuild



Chasing Bottlenecks

- ❑ From hardware to software
- ❑ Hardware getting faster
- ❑ The CPU is still the center player



Compute Intensive Components

- ❑ **Erasure coding**
- ❑ Hashing
- ❑ Encryption
- ❑ Compression



Test Setup – Benchmark of MemoScale vs Jerasure and Intel ISA-L

Data blocks	10
Redundancy blocks	4
Block size	4 MB
Number of cores/threads	1
Measurement criteria	Data encoded excl. redundancies, MB/s
Codes tested	Reed Solomon and MemoScale compute optimized EC
Libraries tested	Jerasure, Intel ISA-L, MemoScale

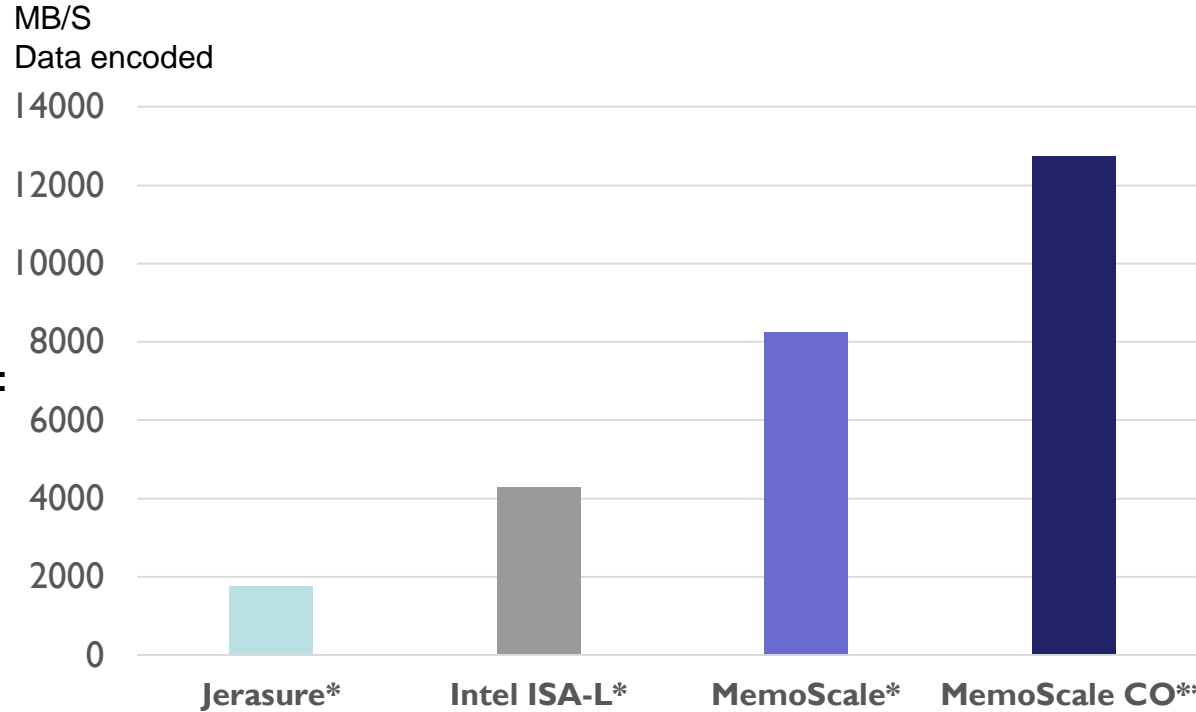


EC Encoding Speed I – Intel Processor

Processor:
Intel Xeon
E5-2676 v3
2.4 Ghz

Erasure coding (EC):
* Reed Solomon

** MemoScale
compute optimized



EC Encoding Speed II – ARM Processor

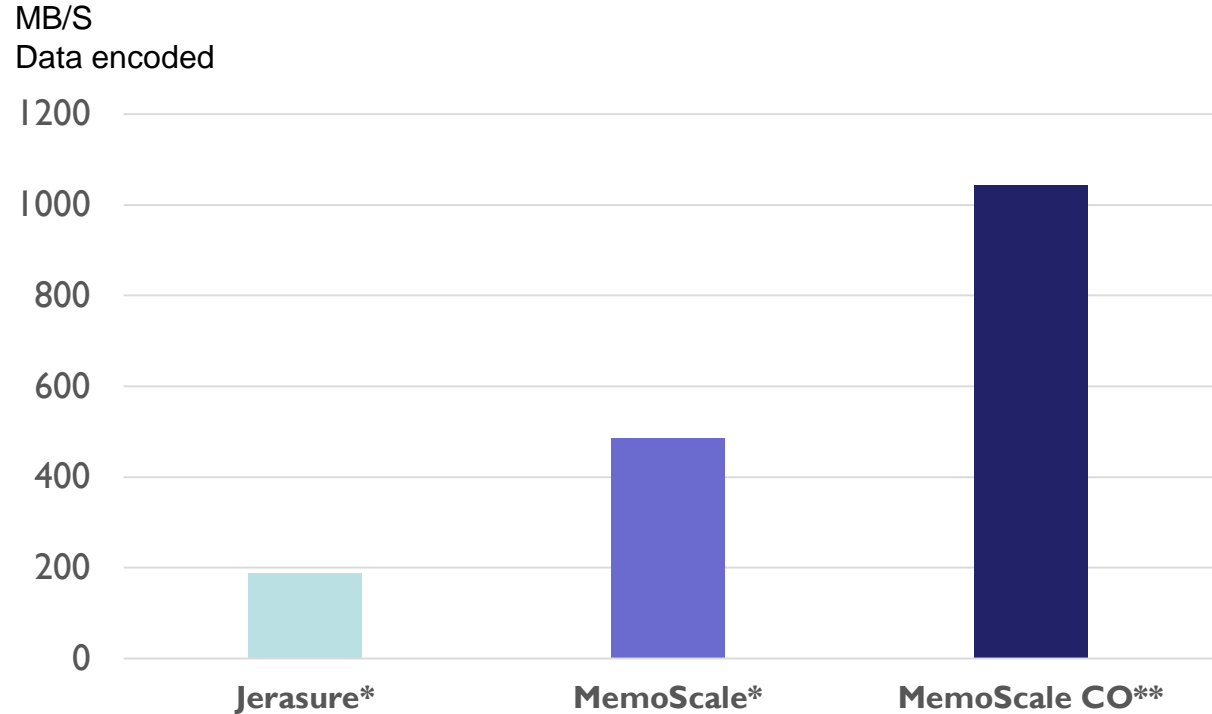
Processor:

HiSilicon Kirin 620
processor, ARM
Cortex-A53 64-bit
SoC 1.2ghz

Erasure coding (EC):

* Reed Solomon

** MemoScale
compute optimized

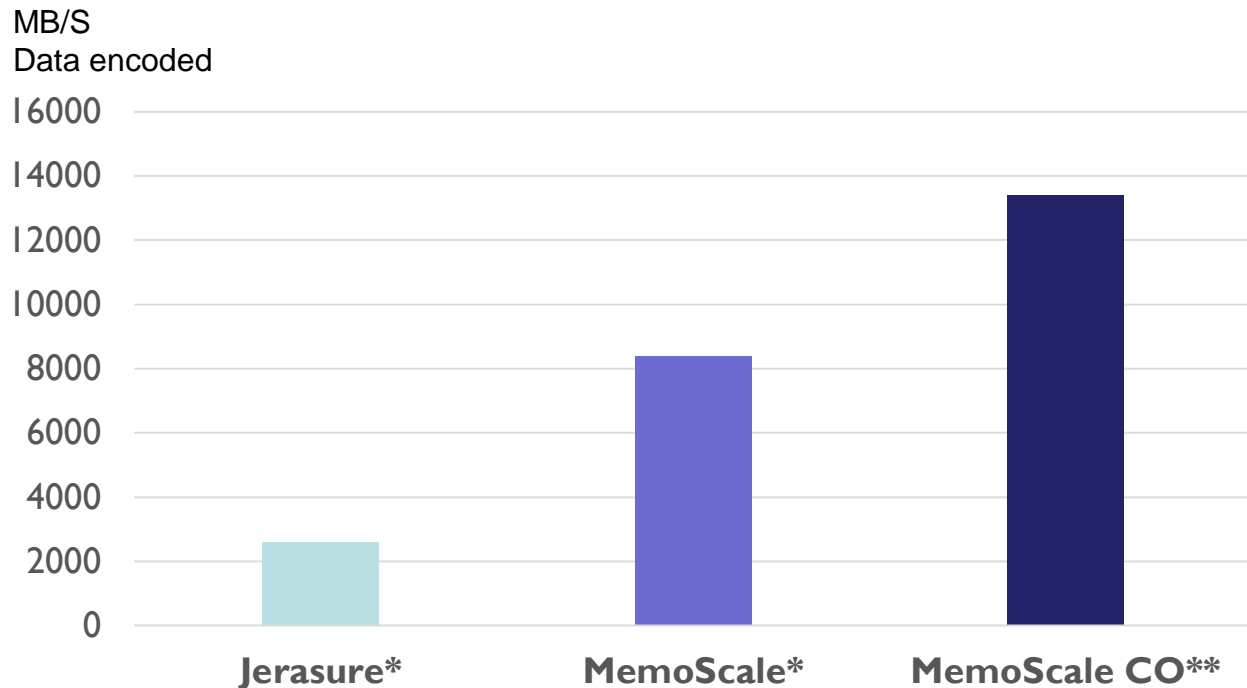


EC Encoding Speed III – AMD Processor

Processor:
AMD RYZEN 7
1700 8-Core 3.0
GHz (3.7 GHz
Turbo)

Erasure coding (EC):
* Reed Solomon

** MemoScale
compute optimized



Finding "Performance Bugs"

- ❑ What is a "Performance Bug"?
- ❑ Optimizing for performance



Measurement Bias I

- ❑ What if changing your user name affected the execution time of programs?



Measurement Bias II

- ❑ What if changing your user name affected the execution time of programs?
 - ❑ Runtime performance change from external factors, introducing bias.
 - ❑ Environment variables
 - ❑ Link ordering



Measurement Bias III

- ❑ What if changing your user name affected the execution time of programs?
 - ❑ Changes in memory layout interacting with various hardware mechanisms.
 - ❑ Hard to predict :(



Measurement Bias IV

- ❑ Challenge for performance analysis
 - ❑ Comparison of algorithms on different system configurations can give misleading results.
 - ❑ Effects can invalidate perceived speedup (Mytkowicz et al. [1]).



Methodology I

- ❑ Control execution
 - ❑ Control execution environment (disable ASLR*).
 - ❑ Instrument program using hardware performance counters.
 - ❑ Measure over range of environments.

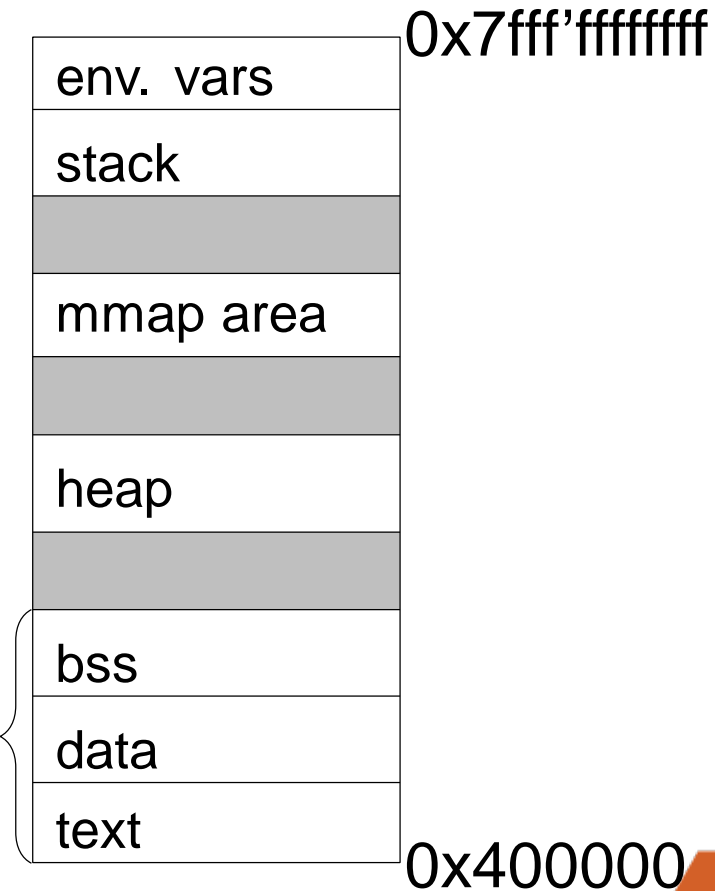
*Address space layout randomization.



Methodology II

- ❑ Memory execution context, assuming a 64-bit process running on a Linux system.

Program code
and static data



Bias from Environment Size I

- A simple program*, with no surprises?

```
static int i, j, k;

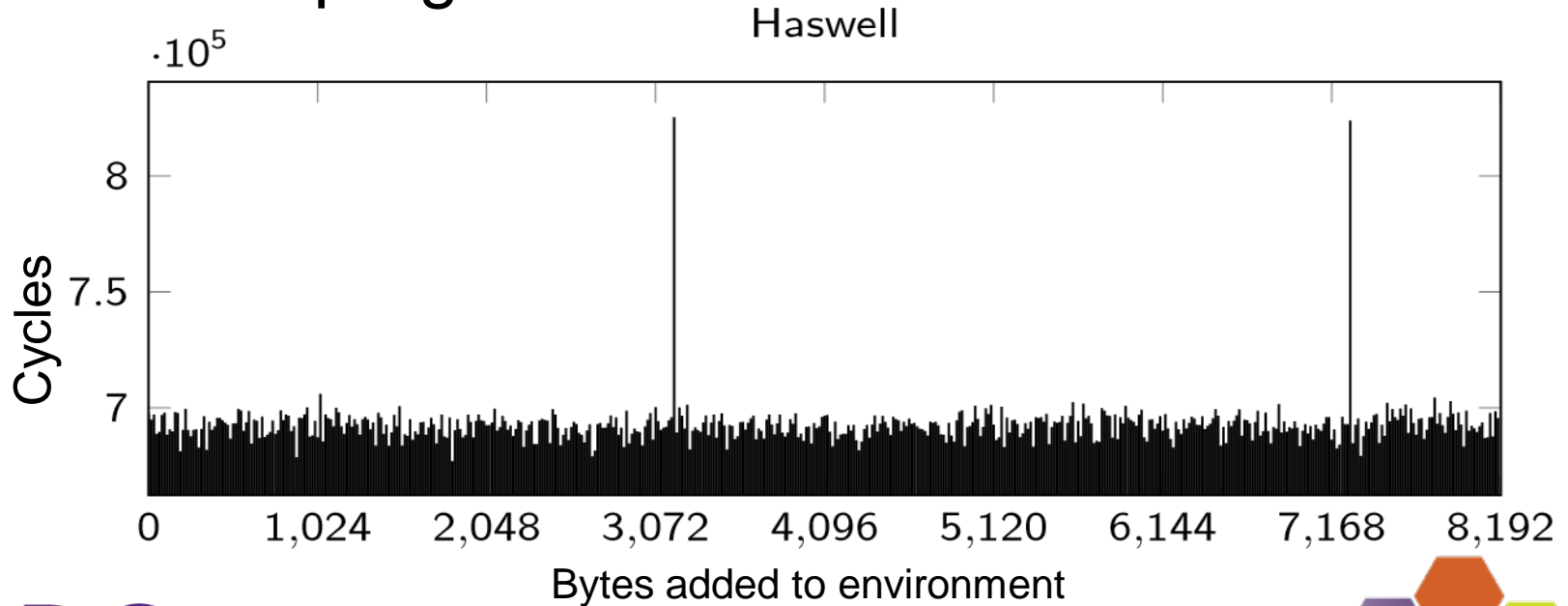
int main() { int g = 0, inc = 1;
  for (; g < 65536; g++) { i +=
    inc; j += inc; k += inc;
  }
  return 0;
}
```

*Microkernel first presented by Mytkowicz *et al.*[2], showing env. bias.



Bias from Environment Size II

- Execute program under different environments.



Bias from Environment Size III

- ❑ How to find the cause?
 - ❑ Accurate measurements
 - ❑ Hardware performance counters
 - ❑ Use correlation
 - ❑ Check *lots* of different CPU performance events
 - ❑ Correlate them with performance



Bias from Environment Size IV

Performance counter	Median	Spike 1	Spike 2
ld blocks partial.address alias	137	327,871	327,848
cpu clk unhalted.thread p	692,686	825,311	823,872
uops executed port.port 0	240,248	106,477	106,409
uops executed port.port 1	249,016	100,913	101,250
uops executed port.port 2	336,768	345,932	346,328
uops executed port.port 3	358,848	356,919	357,723
uops executed port.port 4	280,766	276,036	276,199
uops executed port.port 5	251,631	121,517	121,571
uops executed port.port 6	234,337	178,586	178,646
uops executed port.port 7	137,428	130,325	130,542
resource stalls.any	274,373	406,364	406,256
resource stalls.rs	272,371	135,567	136,503
cycle activity.cycles ldm pending	590,879	718,973	721,580
cycle activity.stalls l2 pending	465,205	490,004	491,656



Bias from Environment Size IV

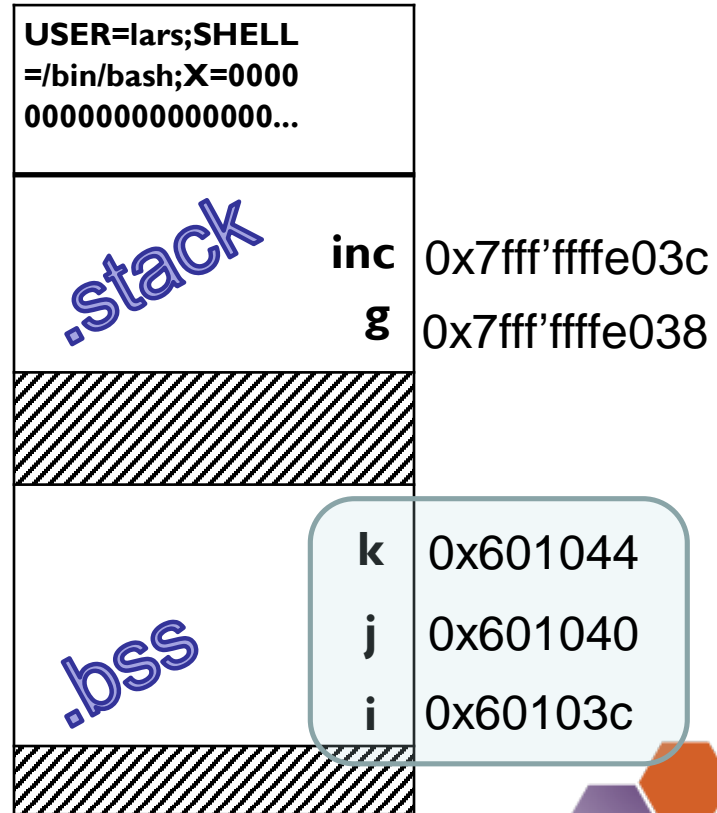
Performance counter	Median	Spike 1	Spike 2
Id blocks partial.address alias	137	327,871	327,848
cpu clk unhalted.thread p	692,686	825,311	823,872
uops executed port.port 0	240,248	106,477	106,409
uops executed port.port 1	249,016	100,913	101,250
uops executed port.port 2	336,768	345,932	346,328
uops executed port.port 3	358,848	356,919	357,723
uops executed port.port 4	280,766	276,036	276,199
uops executed port.port 5	251,631	121,517	121,571
uops executed port.port 6	234,337	178,586	178,646
uops executed port.port 7	137,428	130,325	130,542
resource stalls.any	274,373	406,364	406,256
resource stalls.rs	272,371	135,567	136,503
cycle activity.cycles ldm pending	590,879	718,973	721,580
cycle activity.stalls l2 pending	465,205	490,004	491,656

Correlation!



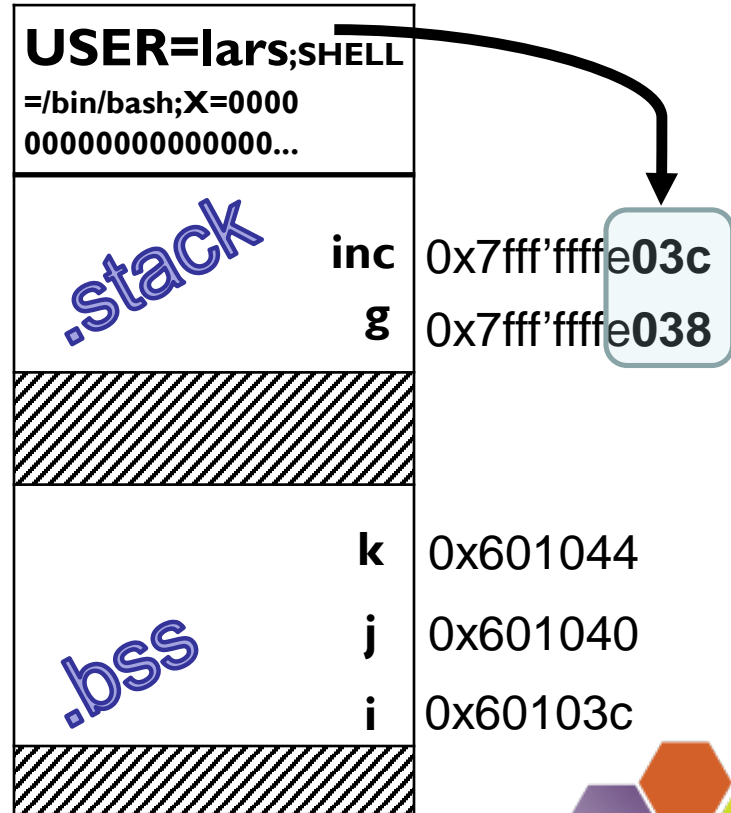
The cause: 4K Address Aliasing I

- Addresses of static i , j and k are fixed for all environments.



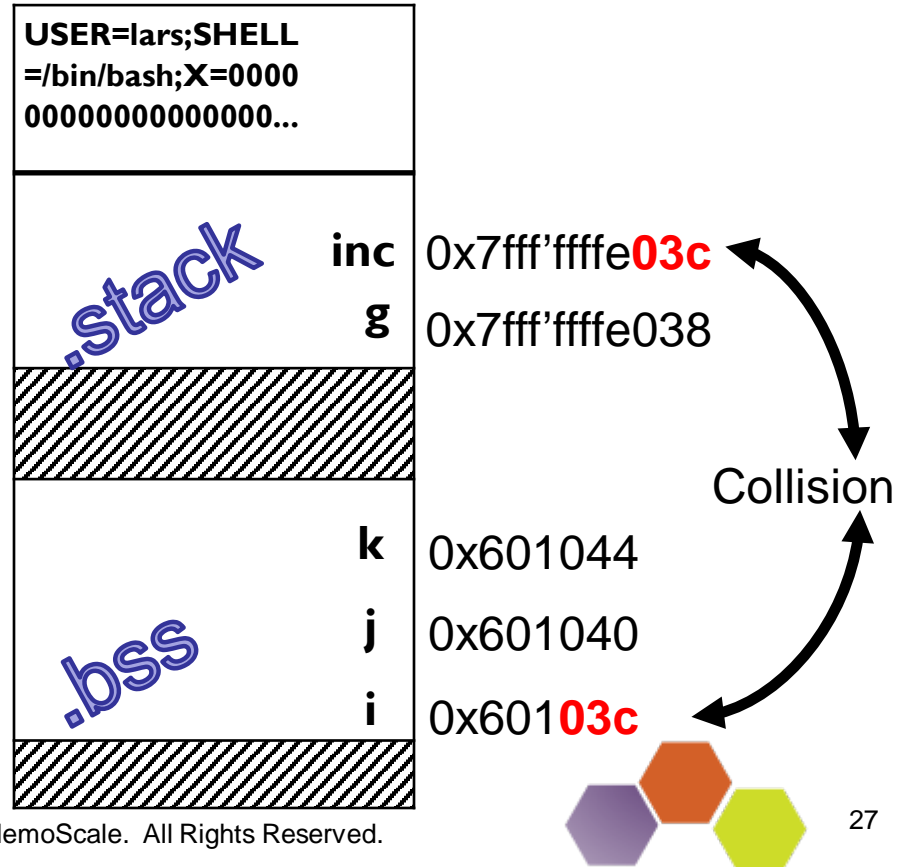
The cause: 4K Address Aliasing II

- Addresses of *inc* and *g* live on the stack.
- The stack is pushed down by the environment variables!



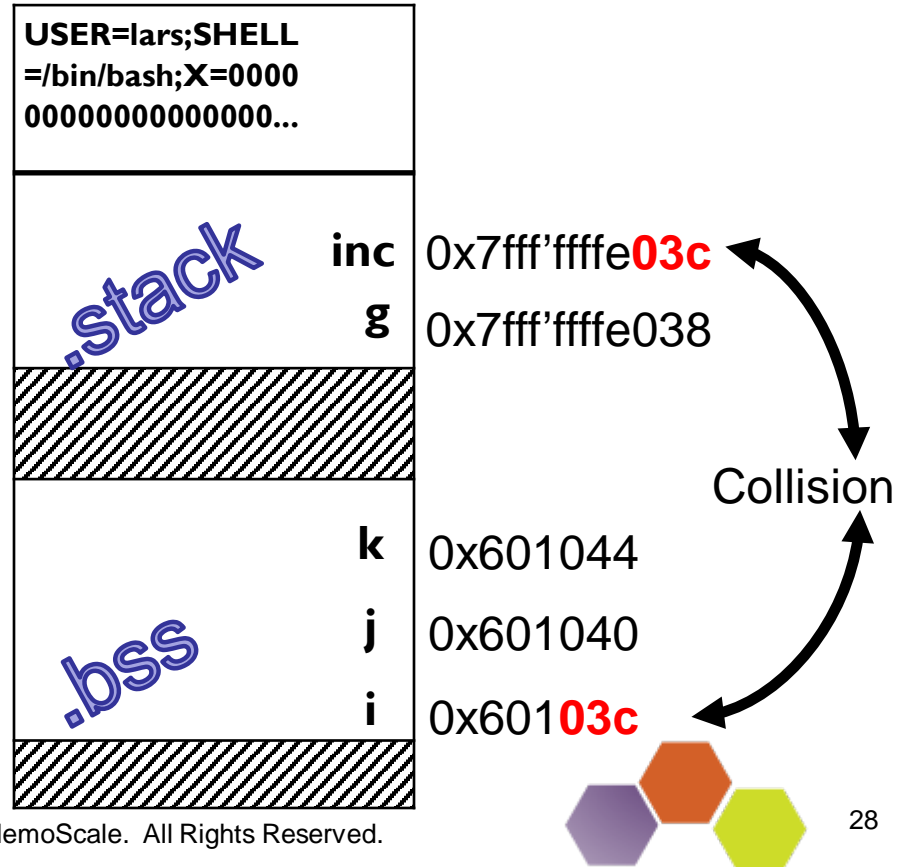
The cause: 4K Address Aliasing III

- ❑ Spike when address of *inc* is aliasing with *i* in the last 12 bits.
 - ❑ The last 3 hex digits.



The cause: 4K Address Aliasing IV

- ❑ Performance impact from *false* dependencies in CPU.
- ❑ Hurts Out-of-Order execution.



How to Measure Address Aliasing

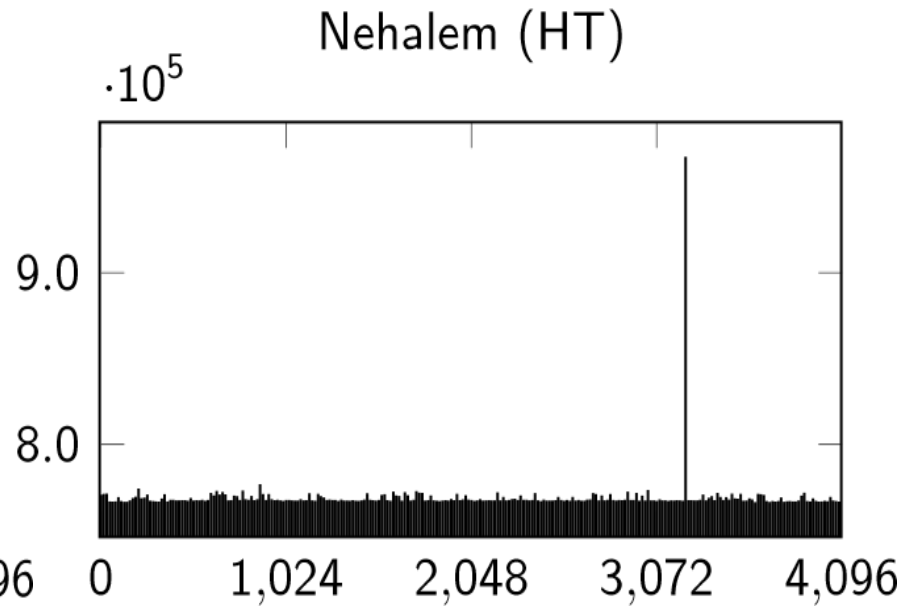
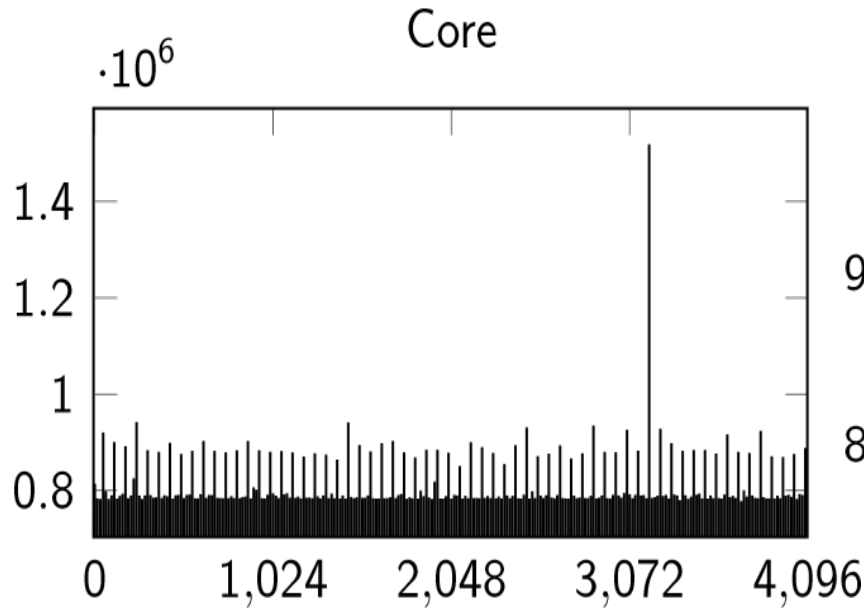
- ❑ **Hardware counter event**
 - ❑ **Linux: perf**
 - ❑ **Manually add event type**

LD BLOCKS PARTIAL.ADDRESS_ALIAS

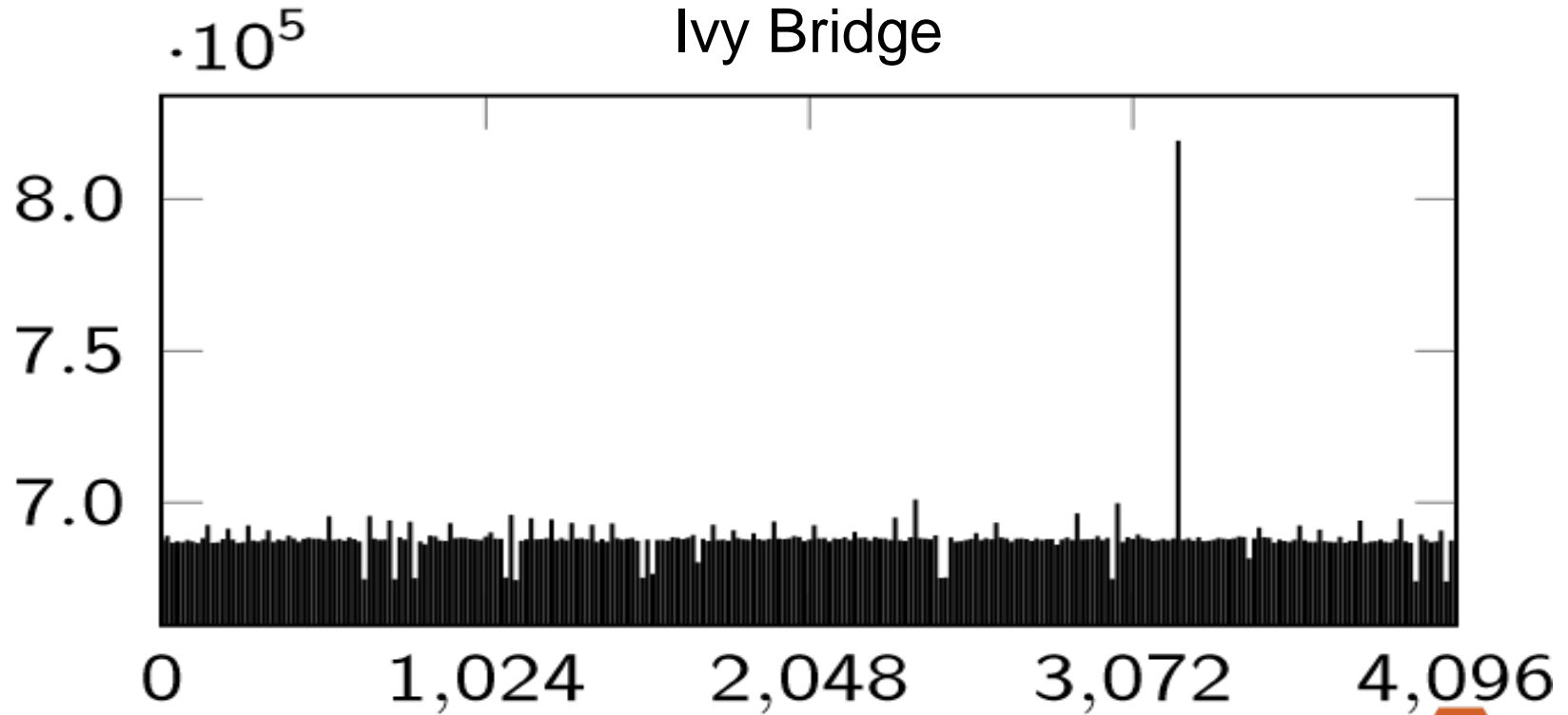
“Counts the number of loads that have partial address match with preceding stores, causing the load to be reissued.” [3, B.3.4.4]



Other architectures I



Other architectures II



Bias from Heap Allocation I

- ❑ Alias conflicts can happen between any sections of memory.
- ❑ While stack addresses depend on environment, heap addresses depend on allocators.



Bias from Heap Allocation II

- Addresses returned by different heap allocators when allocating pairs of equally sized buffers.
- Same 12 bit suffix indicate an aliasing pair.

	5,120 B	1,048,576 B
glibc (ptmalloc)	0x602010	0x2aaaaaf6010
	0x603420	0x2aaaab096010
tcmalloc	0xe0e000	0xe0e000
	0xe10800	0xf0e000
jemalloc	0x2aaaabc0e000	0x2aaaabc06000
	0x2aaaabc10000	0x2aaaabd06000
Hoard	0x2aaaaab00070	0x2aaaaab00070
	0x2aaaaab02070	0x2aaaabf40070



Bias from Heap Allocation III

- ❑ Most allocators tend to use anonymous memory mapping for large requests.
 - ❑ This performed by mmap.
 - ❑ Always at a page boundary.
- ❑ Since page size is 4 KiB, such allocations will *always* alias with each other.



Bias from Heap Allocation IV

- ❑ Other reasons for 4K alignment
 - ❑ RDMA
 - ❑ Requires page locked memory
 - ❑ Naturally aligned to 4K
 - ❑ Drives use 4K blocks (not the same, but...)
 - ❑ Memory mapped files
 - ❑ mmap – always page aligned



Optimization Resistant Code I

- ❑ First example program was too simple?
- ❑ Simplified implementation of convolution with a fixed kernel?



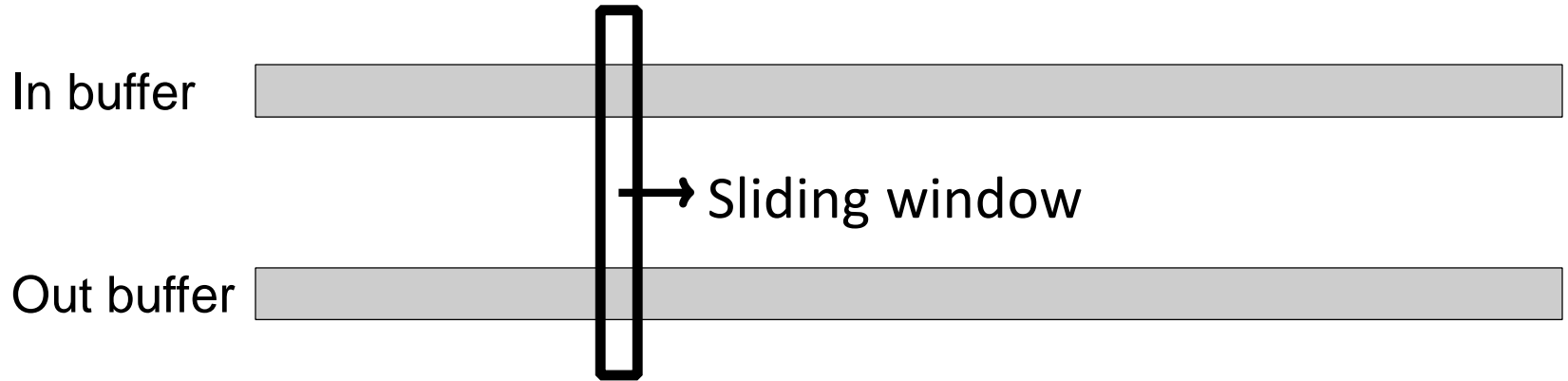
Optimization Resistant Code II

```
static float k[5] = {0.1, 0.25, 0.3, 0.25, 0.1};  
  
void conv(int n, const float *in, float  
  *out) { int i, j;  
  for (i = 2; i < n - 2; ++i) { out[i]  
    = 0; for (j = 0; j < 5; ++j) {  
      out[i] += in[i-2+j] * k[j]; }  
    }  
  }
```



Optimization Resistant Code III

- The key property is the data access pattern:

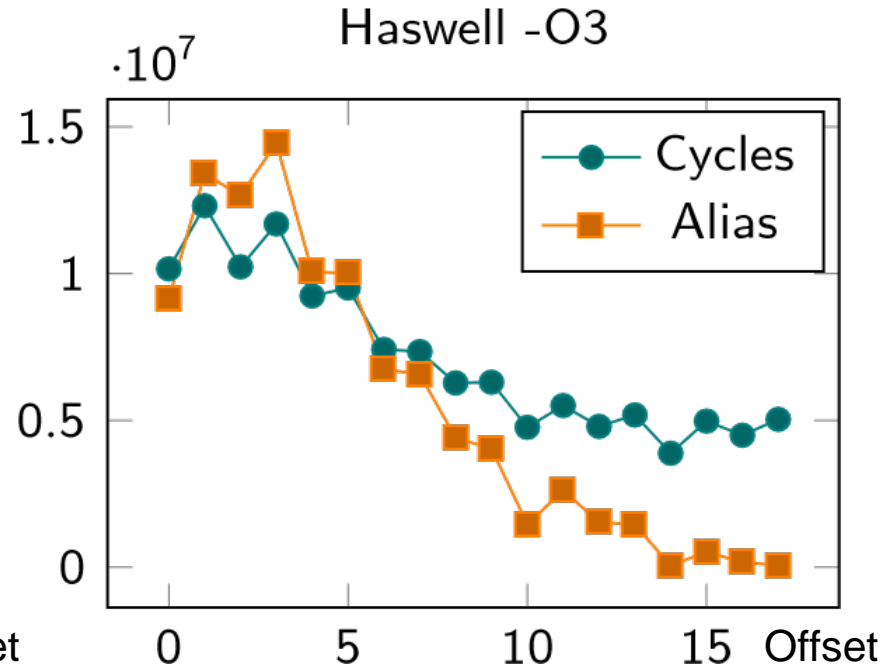
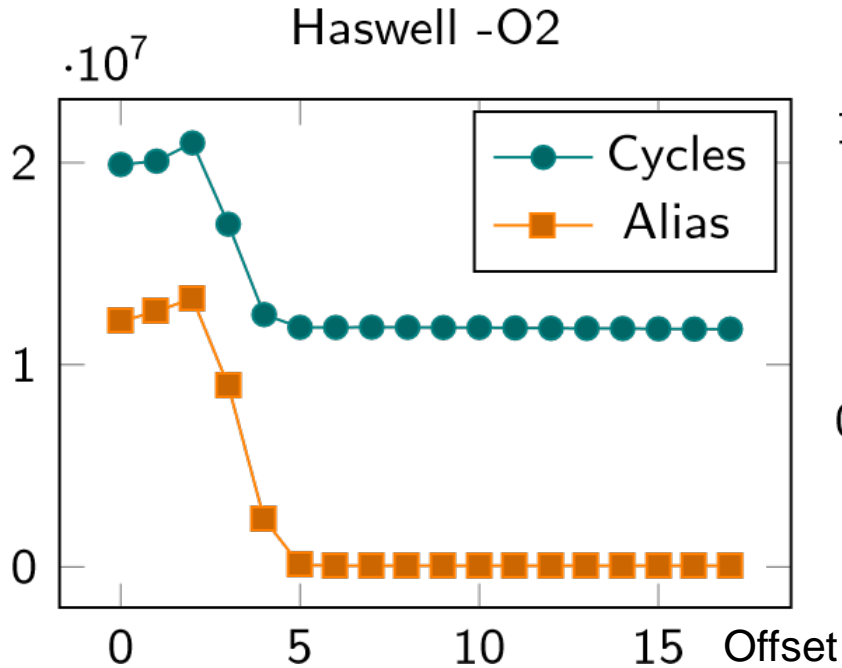


Optimization Resistant Code IV

- ❑ Adjust address of output array manually.
 - ❑ Offset arrays by d elements.
 - ❑ `conv(N, in, out + d);`
 - ❑ $N = 2^{20}$ (4 MiB data)
 - ❑ Evaluate performance.



Optimization Resistant Code V



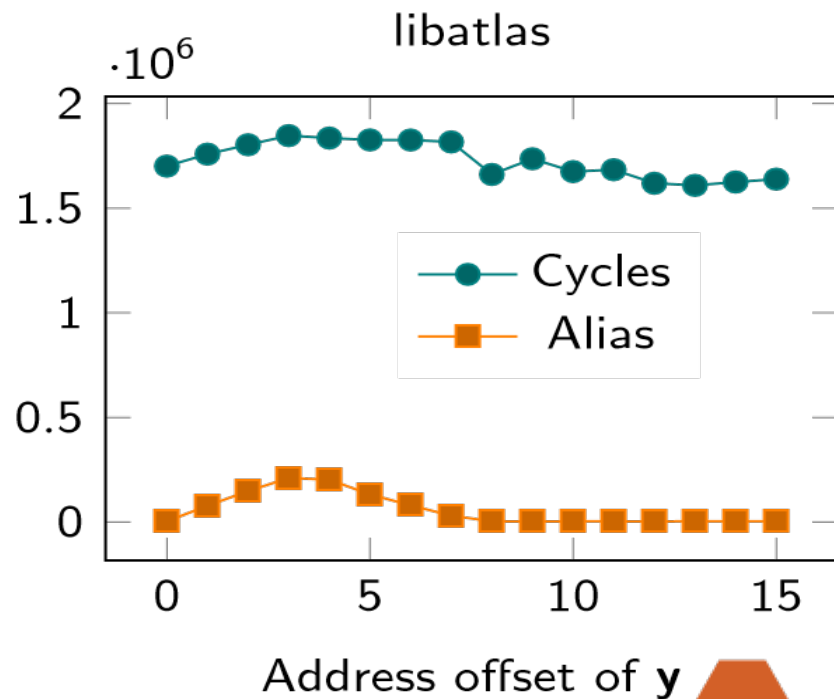
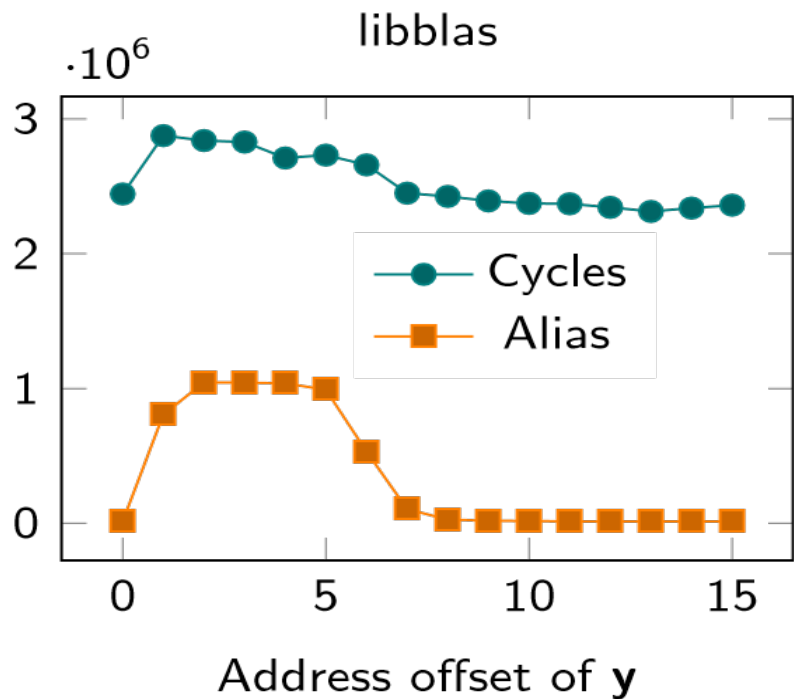
Aliasing in BLAS* libraries I

- ❑ What about heavily optimized libraries?
 - ❑ Measured matrix-vector multiplication (cblas dgemv) with increasing offset between buffer addresses. Computing $\mathbf{y} = \mathbf{Ax}$.

*Basic Linear Algebra Subprograms



Aliasing in BLAS libraries II



Optimizing for address aliasing I

- ❑ Mark buffers with **restrict**
 - ❑ Can allow the optimizer to reduce number of memory accesses.
- ❑ Use optimized libraries
 - ❑ They might mitigate the effect.



Optimizing for address aliasing II

- ❑ Use a special purpose allocator.
 - ❑ Heuristics for avoiding pairwise aliasing buffers.
 - ❑ No known implementation with this goal in mind, to our knowledge 😞



Optimizing for address aliasing III

- ❑ Explicitly adjust address offset.
 - ❑ Consider memory address alignment for performance tuning inner loops.
 - ❑ Buffer addresses as candidate for automatic performance tuning.



Learning points I

- ❑ Accurate benchmarking is hard.
 - ❑ Many uncontrolled effects creates bias.
 - ❑ CPU's are very complex.
- ❑ Optimizing is hard.
 - ❑ Too many rules to address at the same time.
 - ❑ Performance is a sum of good and bad cases.



Learning points II

- ❑ We focused on address aliasing.
 - ❑ And some of the implications it had.
- ❑ There are more (and bigger) issues like this.
 - ❑ CPU cache!
 - ❑ Multithread locks 😞
 - ❑ TLB (Translation lookaside buffer).



Learning points III / Summary

- ❑ Accurate benchmarks are crucial.
 - ❑ Large set of benchmarks.
 - ❑ Many measurement metrics needed.
 - ❑ Not only timings.
 - ❑ Identify performance issues (correlation).
 - ❑ Address aliasing (in our case).



References

1. T. Mytkowicz, P. F. Sweeney, M. Hauswirth, and A. Diwan, “Observer effect and measurement bias in performance analysis,” University of Colorado, Tech. Rep. CU-CS-1042-08, Jun. 2008.
2. T. Mytkowicz, A. Diwan, M. Hauswirth, and P. F. Sweeney, “Producing wrong data without doing anything obviously wrong!” in *Proceedings of the 14th International Conference on Architectural Support for Programming Languages and Operating Systems, ser. ASPLOS XIV*. New York, NY, USA: ACM, 2009, pp. 265–276.
3. Intel © 64 and IA-32 Architectures Optimization Reference Manual, Intel Corporation, Mar. 2014.
4. L. K. Melhus and R. E. Jensen, "Measurement Bias from Address Aliasing," *2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, Chicago, IL, 2016, pp. 1506-1515.



Q&A

- ❑ Further Questions...?
 - ❑ rune.jensen@memoscale.com
 - ❑ per.simonsen@memoscale.com

